# Large Language Models Are State-of-the-Art Evaluators of Translation Quality

## Abstract

We describe GEMBA, a GPT-based metric for assessment of translation quality, which works both with a reference translation and without. In our evaluation, we focus on zero-shot prompting, comparing four prompt variants in two modes, based on the availability of the reference. We investigate seven versions of GPT models, including ChatGPT. We show that our method for translation quality assessment only works with GPT 3.5 and larger models. Comparing to results from WMT22's Metrics shared task, our method achieves state-of-the-art accuracy in both modes when compared to MQM-based human labels. Our results are valid on the system level for all three WMT22 Metrics shared task language pairs, namely English into German, English into Russian, and Chinese into English. This provides a first glimpse into the usefulness of pre-trained, generative large language models for quality assessment of translations. We publicly release all our code and prompt templates used for the experiments described in this work, as well as all corresponding scoring results, to allow for external validation and reproducibility.